

“Descobri uma demonstração maravilhosa desta proposição que, no entanto, não cabe nas margens deste ~~livro~~ apresentação” (Pierre de Fermat, Aritmética de Diophanti).

Convergência

Paulo Ricardo Lisboa de Almeida



Atualização dos pesos

Das aulas passadas, a atualização dos vetores de peso na iteração $\tau + 1$ é dada por.

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau})$$

Onde η é o fator de aprendizagem.

Fator de Aprendizagem

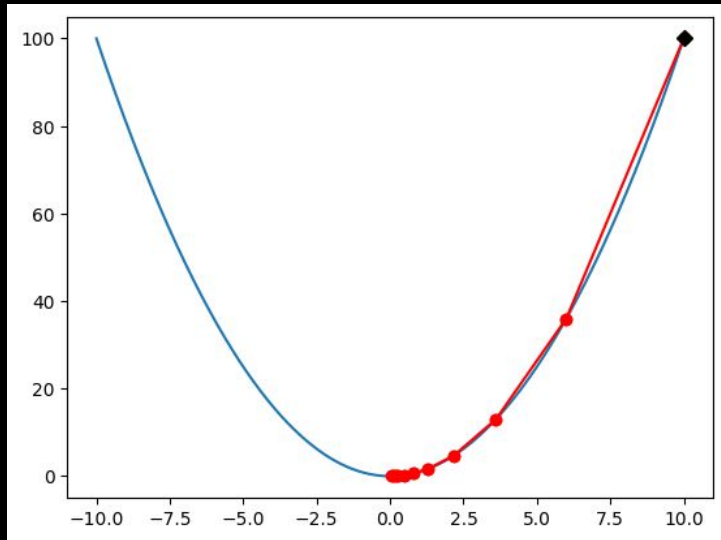
A intuição é que quanto maior o fator de aprendizagem η , mais rápido vamos convergir para a resposta.

Teste essa intuição com o programa disponibilizado no Google Colab.

Fator de Aprendizagem

Encontrando o mínimo da função x^2 considerando $x=10$ como chute inicial, e:

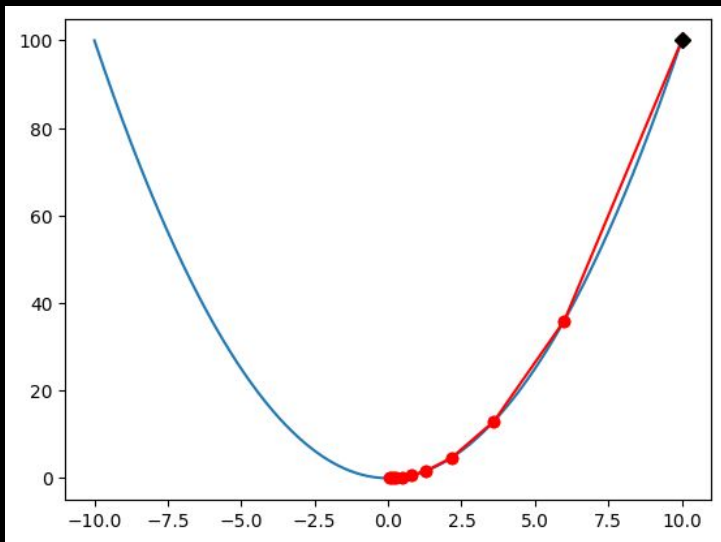
$$\eta = 0.2$$



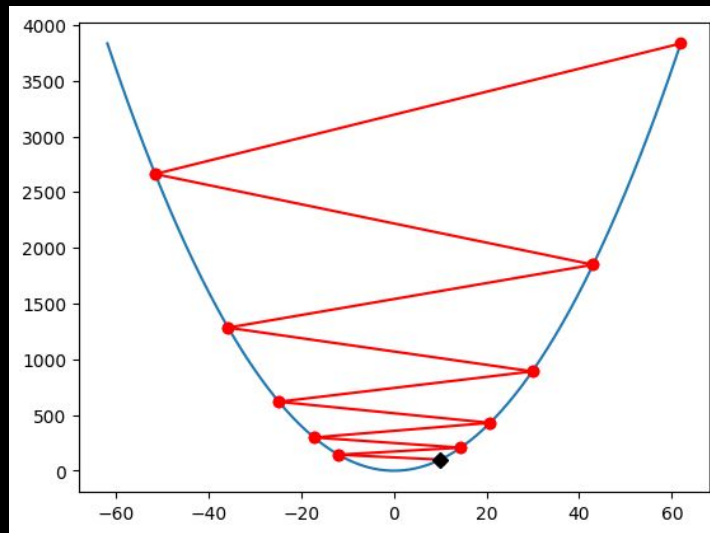
Fator de Aprendizagem

Encontrando o mínimo da função x^2 considerando $x=10$ como chute inicial, e:

$\eta = 0.2$



$\eta = 1.1$



Fator de Aprendizagem

A derivada de segunda ordem do *loss* em função dos pesos dá o fator de aprendizagem ótimo. Assumindo que a função é aproximadamente quadrática:

$$\eta_{opt} = \left(\frac{\partial^2 L}{\partial w^2} \right)^{-1}$$

Fator de Aprendizagem

A derivada de segunda ordem do *loss* em função dos pesos dá o fator de aprendizagem ótimo. Assumindo que a função é aproximadamente quadrática:

$$\eta_{opt} = \left(\frac{\partial^2 L}{\partial w^2} \right)^{-1}$$

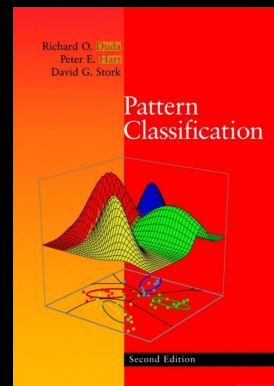
Se $0 < \eta \leq \eta_{opt}$ a convergência é garantida.

Mas pode ser lenta, se $\eta < \eta_{opt}$.

Se $\eta_{opt} < \eta < 2\eta_{opt}$ o sistema oscila.

Se $\eta > 2\eta_{opt}$ o sistema diverge.

Duda, R. O., Hart, P. E., Stork, D.
G. Pattern Classification. 2012.



Exemplo

Para a função $f(x) = x^2$ do exemplo:

$$f(x) = x^2$$

$$f'(x) = 2x dx$$

$$f''(x) = 2 dx$$

Logo:

$$\eta_{opt} = (2)^{-1} = \frac{1}{2}$$

Exemplo

Para a função $f(x) = x^2$ do exemplo:

$$f(x) = x^2$$

$$f'(x) = 2x dx$$

$$f''(x) = 2 dx$$

Logo:

$$\eta_{opt} = (2)^{-1} = \frac{1}{2}$$

Valores acima de 0.5 vão fazer o sistema oscilar, e acima de 1 vão fazer o sistema divergir.

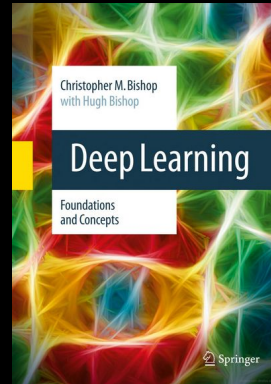
Fator de Aprendizagem

Na prática, apenas garantimos que η seja pequeno o suficiente para garantir convergência.

O cálculo do valor ótimo não é trivial em altas dimensões.

Depende dos autovalores das matrizes.

Veja uma discussão em Bishop, Bishop (2023).



Momentum

Ideia: adicionar o conceito de inércia a descida de gradiente.

Se estávamos “descendo” em determinada direção, tendemos a continuar nessa direção na próxima iteração.

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau}) + \mu \Delta \mathbf{w}^{\tau-1}$$

Onde $0 \leq \mu < 1$.

Momentum

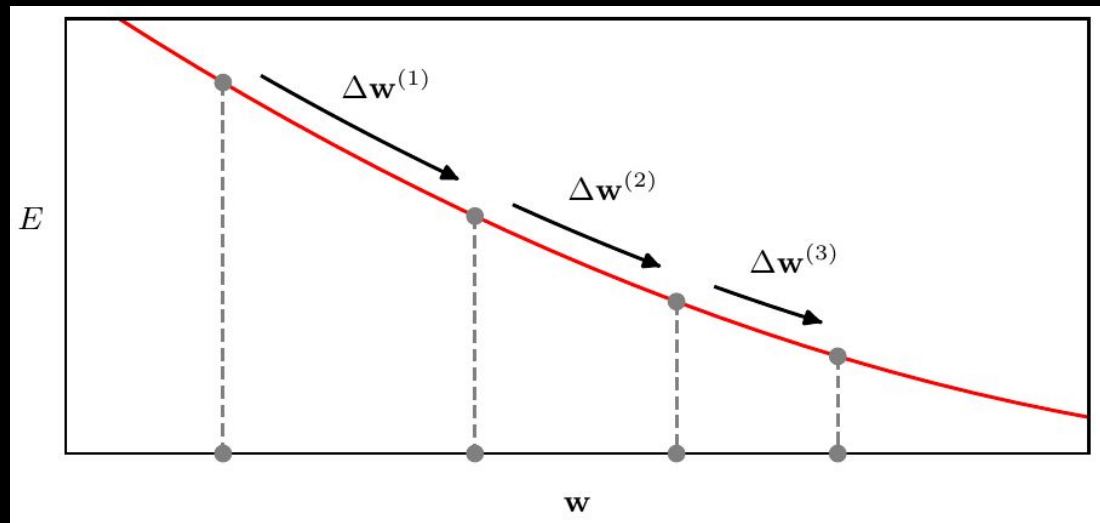
$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau}) + \mu \Delta \mathbf{w}^{\tau-1}$$

Gradiente atual.

Atualização do vetor da etapa anterior.

Momentum - Exemplo

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau}) + \mu \Delta \mathbf{w}^{\tau-1}$$



Considere:

- Uma superfície de erro com curvatura relativamente baixa.
- O gradiente ∇E é aproximadamente constante a cada iteração.

Exemplo de Bishop e Bishop (2024). Na figura **não** é aplicado momentum.

Momentum - Exemplo

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau}) + \mu \Delta \mathbf{w}^{\tau-1}$$

Considere:

- Uma superfície de erro com curvatura relativamente baixa.
- O gradiente ∇E é aproximadamente constante a cada iteração.

Momentum - Exemplo

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau}) + \mu \Delta \mathbf{w}^{\tau-1}$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E + \mu(-\eta \nabla E + \mu \Delta \mathbf{w}^{\tau-2})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E - \mu \eta \nabla E + \mu^2 \Delta \mathbf{w}^{\tau-2}$$

Considere:

- Uma superfície de erro com curvatura relativamente baixa.
- O gradiente ∇E é aproximadamente constante a cada iteração.

Momentum - Exemplo

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau}) + \mu \Delta \mathbf{w}^{\tau-1}$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E + \mu(-\eta \nabla E + \mu \Delta \mathbf{w}^{\tau-2})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E - \mu \eta \nabla E + \mu^2 \Delta \mathbf{w}^{\tau-2}$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E - \mu \eta \nabla E + \mu^2(-\eta \nabla E + \mu \mathbf{w}^{\tau-3})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E - \mu \eta \nabla E - \mu^2 \eta \nabla E + \mu^3 \mathbf{w}^{\tau-3}$$

...

Considere:

- Uma superfície de erro com curvatura relativamente baixa.
- O gradiente ∇E é aproximadamente constante a cada iteração.

Momentum - Exemplo

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau}) + \mu \Delta \mathbf{w}^{\tau-1}$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E + \mu(-\eta \nabla E + \mu \Delta \mathbf{w}^{\tau-2})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E - \mu \eta \nabla E + \mu^2 \Delta \mathbf{w}^{\tau-2}$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E - \mu \eta \nabla E + \mu^2(-\eta \nabla E + \mu \mathbf{w}^{\tau-3})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E - \mu \eta \nabla E - \mu^2 \eta \nabla E + \mu^3 \mathbf{w}^{\tau-3}$$

...

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(1 + \mu + \mu^2 + \mu^3 + \dots)$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \frac{\eta}{1 - \mu} \nabla E$$

Considere:

- Uma superfície de erro com curvatura relativamente baixa.
- O gradiente ∇E é aproximadamente constante a cada iteração.

Momentum - Exemplo

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau}) + \mu \Delta \mathbf{w}^{\tau-1}$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E + \mu(-\eta \nabla E + \mu \Delta \mathbf{w}^{\tau-2})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E - \mu \eta \nabla E + \mu^2 \Delta \mathbf{w}^{\tau-2}$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E - \mu \eta \nabla E + \mu^2(-\eta \nabla E + \mu \mathbf{w}^{\tau-3})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E - \mu \eta \nabla E - \mu^2 \eta \nabla E + \mu^3 \mathbf{w}^{\tau-3}$$

...

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(1 + \mu + \mu^2 + \mu^3 + \dots)$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \frac{\eta}{1 - \mu} \nabla E$$

Soma de uma série geométrica. Veja um esboço da prova nos últimos slides.

Considere:

- Uma superfície de erro com curvatura relativamente baixa.
- O gradiente ∇E é aproximadamente constante a cada iteração.

Momentum - Exemplo

Com momentum.

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \frac{\eta}{1 - \mu} \nabla E$$

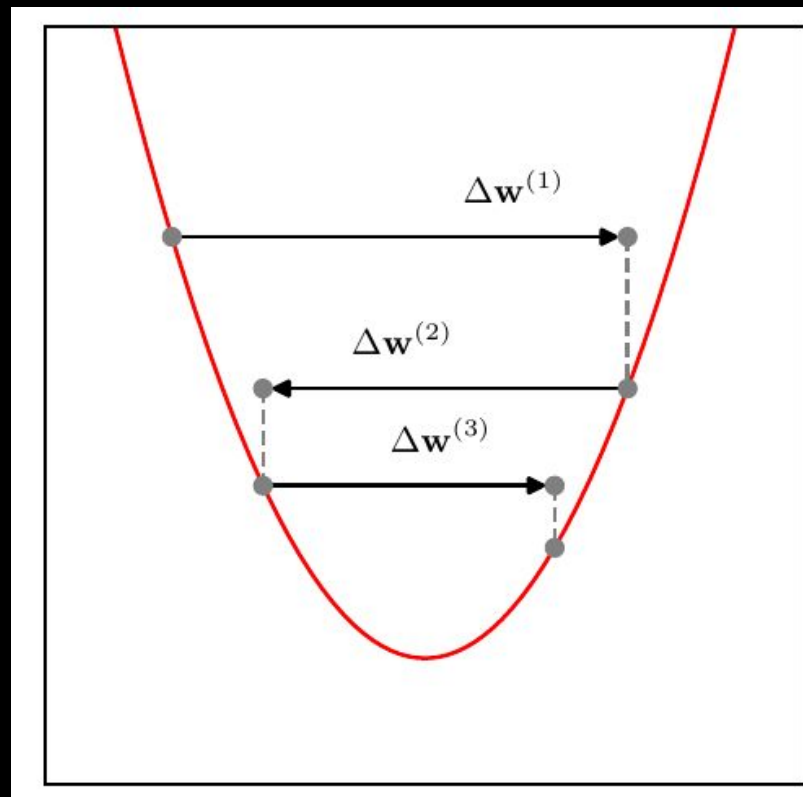
Aceleramos a convergência de η para $\frac{\eta}{1 - \mu}$.

Sem momentum.

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E$$

Momentum Exemplo 2

E para um caso onde há oscilação, o que esperamos que acontecerá ao adicionar momentum?



Exemplo de Bishop e Bishop (2024).

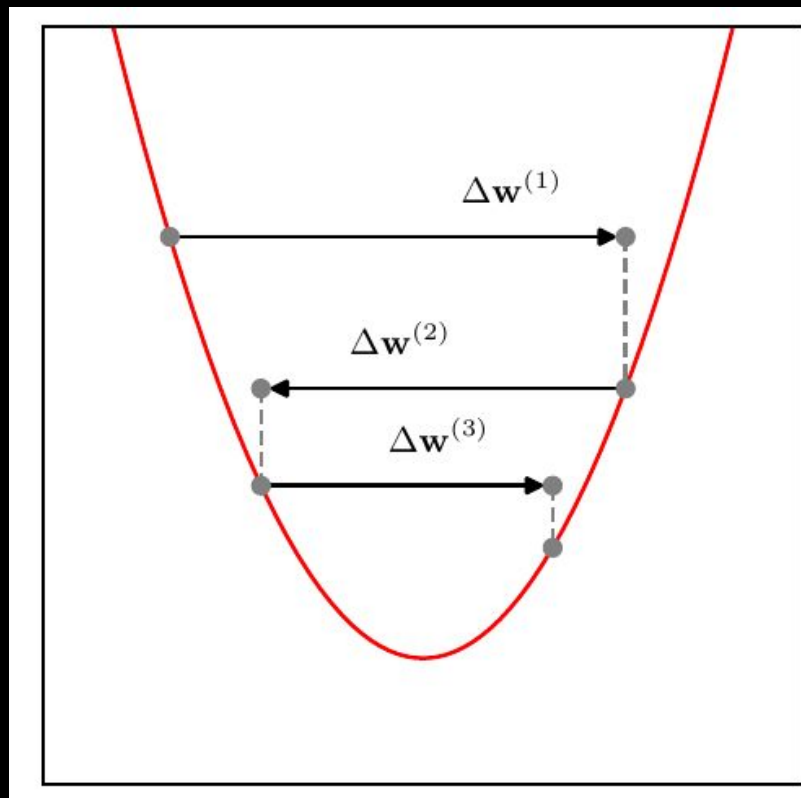
Momentum Exemplo 2

E para um caso onde há oscilação, o que esperamos que acontecerá ao adicionar momentum?

Os termos do momentum tendem a se cancelar.

Hoje temos $+\nabla E$, e em outras temos $-\nabla E$.

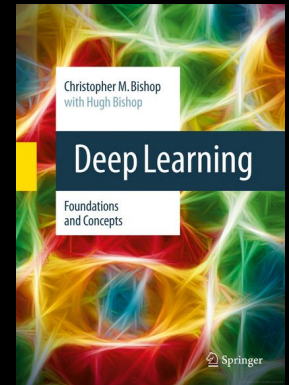
Logo, o *learning rate* será aproximadamente o mesmo que se não tivéssemos momentum: η



Exemplo de Bishop e Bishop (2024).

Na prática

Na prática, um valor comum é $\mu = 0.9$.



Bishop, Bishop (2024).

Learning rate scheduler

Intuição: no começo do treinamento, um *learning rate* maior pode levar a uma convergência mais rápida, e fugir de mínimos locais.

Ao custo de causar oscilações próximo dos pontos de mínimo, onde um *learning rate* pequeno possibilitaria um melhor ajuste fino.

Learning rate scheduler

Intuição: no começo do treinamento, um *learning rate* maior pode levar a uma convergência mais rápida, e fugir de mínimos locais.

Ao custo de causar oscilações próximo dos pontos de mínimo, onde um *learning rate* pequeno possibilitaria um melhor ajuste fino.

Para tentar alcançar o melhor dos dois mundos, podemos usar um *Learning Rate Scheduler* (Escalonador de Fator de Aprendizado).

Escalonadores

Alguns exemplos de escalonadores:

- Linear.
- Lei da potência.
- Decaimento exponencial.

Exemplo - Escalonador Linear

$$\eta^\tau = \begin{cases} (1 - \frac{\tau}{K})\eta^0 + \frac{\tau}{K}\eta^K & , se \tau \leq K \\ \eta^K & , se \tau > K \end{cases}$$

O learning rate inicia em η^0 e decai linearmente até η^K durante K passos. Após isso, é mantido fixo em η^K .

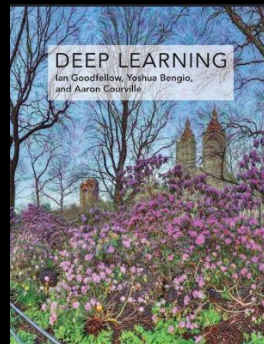
Exemplo - Escalonador Linear

É difícil escolher os valores de K , η^0 e η^K .

Segundo Goodfellow, Bengio e Courville, $\eta^K = 0.01\eta^0$ pode ser um bom chute inicial.

De qualquer forma, analisar empiricamente a curva de aprendizado pode nos indicar bons valores para essas constantes.

Goodfellow, I., Bengio, Y.,
Courville, A. Deep Learning.
2016.



$$\eta^\tau = \begin{cases} (1 - \frac{\tau}{K})\eta^0 + \frac{\tau}{K}\eta^K & , se \tau \leq K \\ \eta^K & , se \tau > K \end{cases}$$

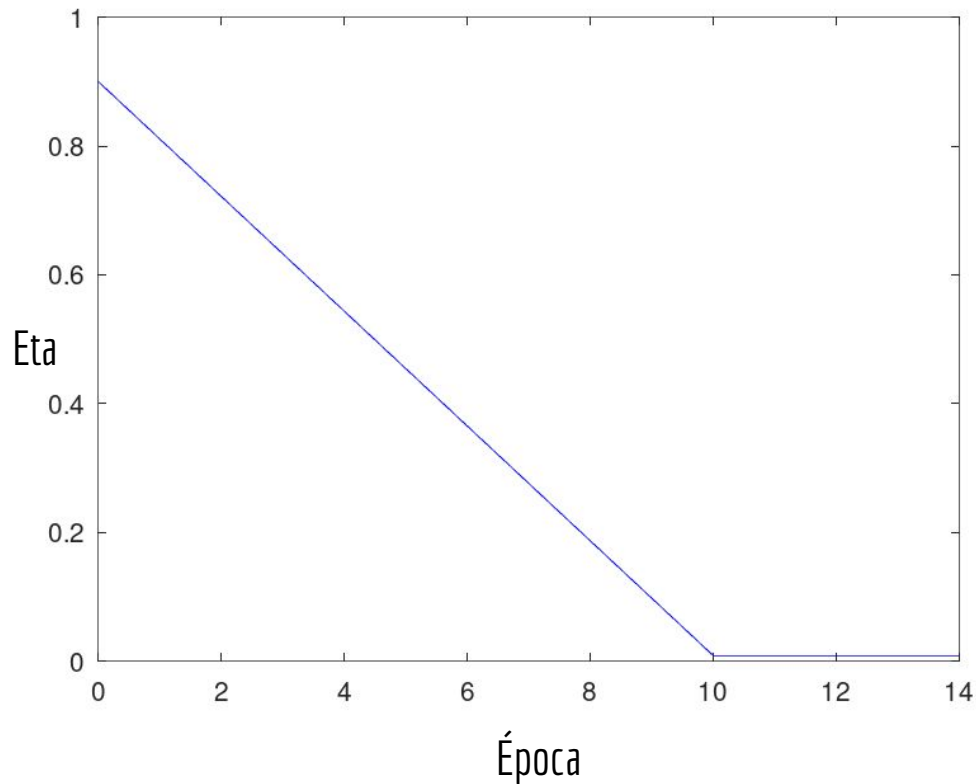
Exemplo - Escalonador Linear

Exemplo considerando:

$$\eta^0 = 0.9$$

$$\eta^K = 0.009$$

$$K = 10$$

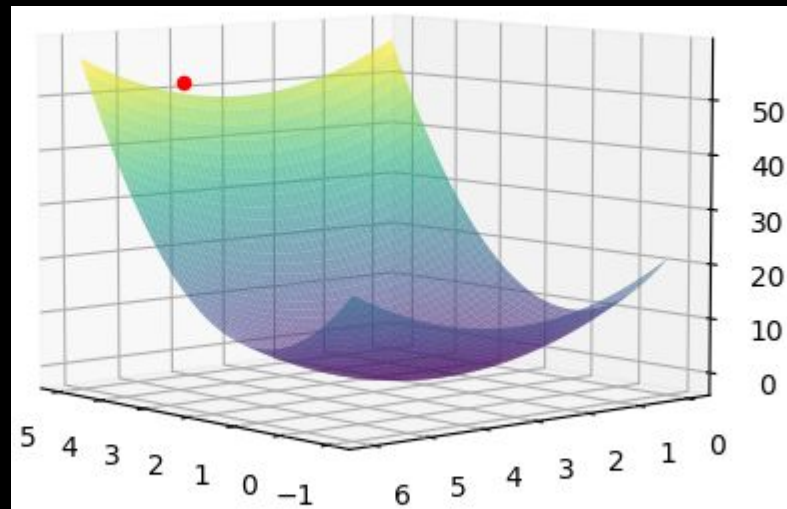


RMSPprop

O momentum tenta otimizar o *learning rate* de forma global.

Mas a curvatura do erro geralmente varia dependendo da direção (ex.: entre o eixo x e y).

Ideia: ter um *learning rate diferente* para cada parâmetro.



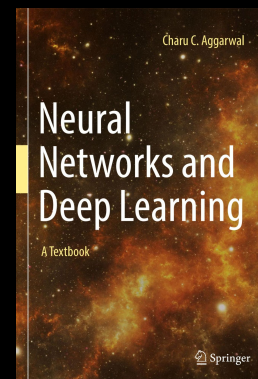
RMSProp

O RMSProp usa uma ideia similar ao Momentum, com as diferenças:

É calculado para cada um dos parâmetros.

Faz a raiz da soma dos quadrados dos gradientes (para detalhes, pesquise sobre o AdaGrad).

Aggarwal. Neural Networks and
Deep Learning. 2018.



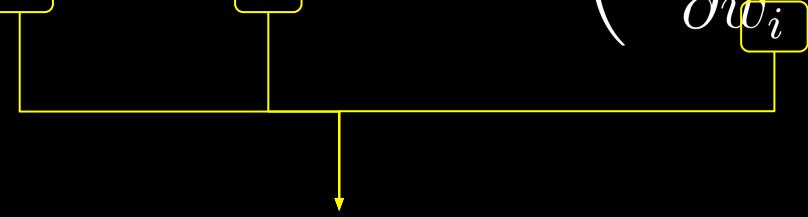
RMSProp

$$r_i^{\tau+1} = \beta r_i^{\tau} + (1 - \beta) \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)^2$$

$$0 < \beta < 1$$

$$w_i^{\tau+1} = w_i^{\tau} - \frac{\eta}{\sqrt{r_i^{\tau+1} + \delta}} \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)$$

RMSPProp

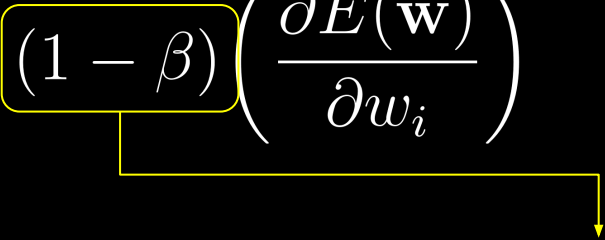
$$r_i^{\tau+1} = \beta r_i^{\tau} + (1 - \beta) \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)^2$$


$$0 < \beta < 1$$

Calculado para o i -ésimo parâmetro.

$$w_i^{\tau+1} = w_i^{\tau} - \frac{\eta}{\sqrt{r_i^{\tau+1} + \delta}} \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)$$

RMSPProp

$$r_i^{\tau+1} = \beta r_i^{\tau} + (1 - \beta) \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)^2 \quad 0 < \beta < 1$$


O último gradiente terá peso $(1 - \beta)$, o penúltimo $(1 - \beta)^2$, ... decaindo exponencialmente.

$$w_i^{\tau+1} = w_i^{\tau} - \frac{\eta}{\sqrt{r_i^{\tau+1} + \delta}} \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)$$

RMSPProp

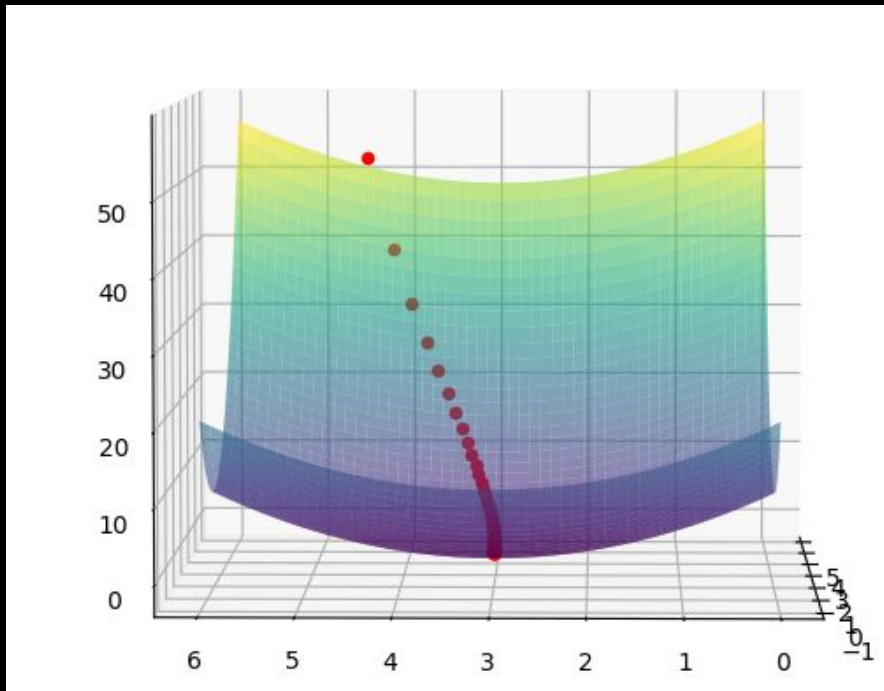
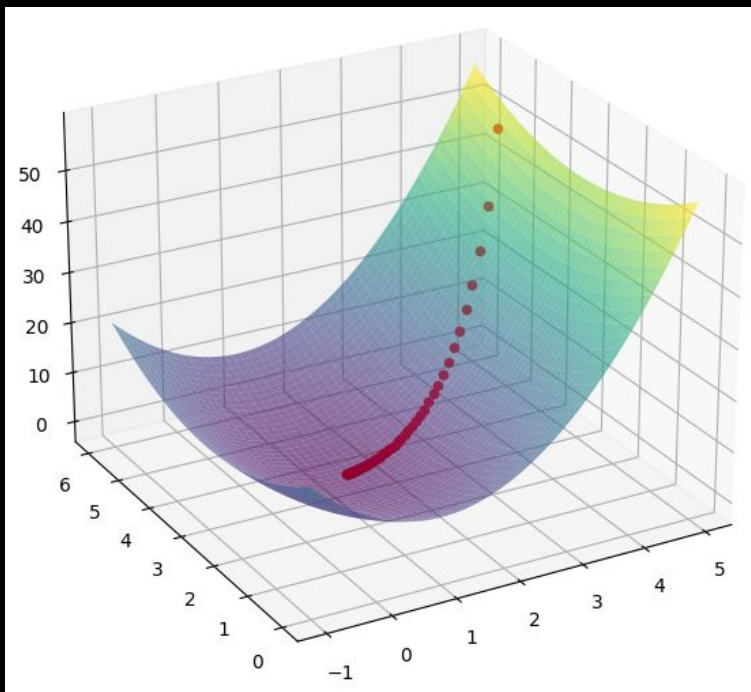
$$r_i^{\tau+1} = \beta r_i^{\tau} + (1 - \beta) \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)^2 \quad 0 < \beta < 1$$

$$w_i^{\tau+1} = w_i^{\tau} - \frac{\eta}{\sqrt{r_i^{\tau+1} + \delta}} \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)$$

Constante para garantir estabilidade numérica caso $r_i^{\tau+1}$ fique próximo de zero. Geralmente $\delta = 10^{-8}$.

Faça você mesmo

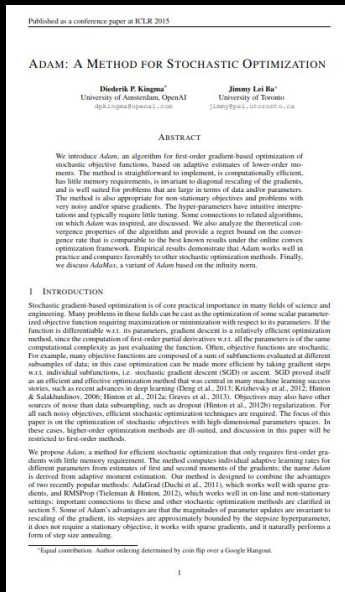
Entenda e execute o exemplo do RMSProp disponibilizado no Google Colab.



Adam

O algoritmo Adam é um dos mais famosos e usados para ajustar o learning rate.

Utiliza uma combinação de RMSProp e Momentum.



Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." *arXiv preprint arXiv:1412.6980* (2014).

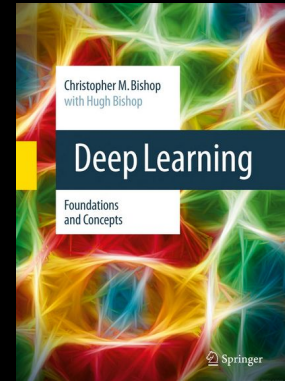
Adam

$$s_i^{\tau+1} = \beta_1 s_i^{\tau} + (1 - \beta_1) \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)$$

$$r_i^{\tau+1} = \beta_2 r_i^{\tau} + (1 - \beta_2) \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)^2$$

$$w_i^{\tau+1} = w_i^{\tau} - \eta \frac{s_i^{\tau+1}}{\sqrt{r_i^{\tau+1} + \delta}}$$

Aqui foi omitida a correção do Bias.
Veja em Bishop e Bishop (2024).



Adam

$$s_i^{\tau+1} = \beta_1 s_i^{\tau} + (1 - \beta_1) \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right) \leftarrow \text{Momentum, com decaimento controlado por } \beta_1.$$

$$r_i^{\tau+1} = \beta_2 r_i^{\tau} + (1 - \beta_2) \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)^2 \leftarrow \text{Momentum, com decaimento controlado por } \beta_2.$$

$$w_i^{\tau+1} = w_i^{\tau} - \eta \frac{s_i^{\tau+1}}{\sqrt{r_i^{\tau+1} + \delta}}$$

Adam

$$s_i^{\tau+1} = \beta_1 s_i^{\tau} + (1 - \beta_1) \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)$$

$$r_i^{\tau+1} = \beta_2 r_i^{\tau} + (1 - \beta_2) \left(\frac{\partial E(\mathbf{w})}{\partial w_i} \right)^2$$

$$w_i^{\tau+1} = w_i^{\tau} - \eta \frac{s_i^{\tau+1}}{\sqrt{r_i^{\tau+1} + \delta}}$$

No artigo original é sugerido usar

$\beta_1 = 0.9$ e $\beta_2 = 0.99$.

Published as a conference paper at ICLR 2015

ADAM: A METHOD FOR STOCHASTIC OPTIMIZATION

Diederik P. Kingma¹
University of Amsterdam, OpenAI
dpk1@openai.com

Jimmy Ba²
University of Toronto
jba@pslab.utoronto.ca

ABSTRACT

We introduce Adam, an algorithm for first-order gradient-based optimization of stochastic objective functions, based on adaptive estimates of lower-order moments. The method is straightforward to implement, is computationally efficient, has little memory requirements, is invariant to diagonal rescaling of the gradients, and is well suited for problems that are large in terms of data and/or parameters. The method is also appropriate for non-stationary objectives and problems with very noisy and/or sparse gradients. The hyper-parameters have intuitive interpretations and typically require little tuning. Some connections to related algorithms, on which Adam was inspired, are discussed. We also analyze the theoretical convergence properties of the algorithm and provide a regret bound on the convergence rate that is comparable to the best known results under the online convex optimization framework. Empirical results demonstrate that Adam works well in practice and compares favorably to other stochastic optimization methods. Finally, we discuss Adam's a variant of Adam based on the infinity norm.

1 INTRODUCTION

Stochastic gradient-based optimization is of core practical importance in many fields of science and engineering. Many problems in these fields can be cast as the optimization of some scalar parameterized objective function requiring maximization or minimization with respect to its parameters. If the function is differentiable w.r.t. its parameters, gradient descent is a relatively efficient optimization method, since the computation of first-order partial derivatives w.r.t. all the parameters is of the same computational complexity as just evaluating the function. Often, objective functions are stochastic. For example, many objective functions are composed of a sum of subfunctions evaluated on different subamples of data; in this case optimization can be made more efficient by taking gradient steps w.r.t. individual subfunctions, i.e. stochastic gradient descent (SGD) or variants, SGD with momentum as an efficient and effective optimization method that was central in many machine learning success stories, such as recent advances in deep learning (Ding et al., 2013; Krizhevsky et al., 2012; Hinton & Salakhutdinov, 2006; Hinton et al., 2012a; Graves et al., 2013). Objectives may also have other sources of noise than data subsampling, such as dropout (Hinton et al., 2012b) or regularization. For all such noisy objectives, efficient stochastic optimization techniques are required. The focus of this paper is on the optimization of stochastic objectives with high-dimensional parameters spaces. In these cases, higher-order optimization methods are ill-suited, and discussion in this paper will be restricted to first-order methods.

We propose Adam, a method for efficient stochastic optimization that only requires first-order gradients with little memory requirement. The method computes individual adaptive learning rates for different parameters from estimates of first and second moments of the gradients; the same Adam is derived from adaptive moment estimation. Our method is designed to combine the advantages of two recently popular methods, AdaGrad (Duchi et al., 2011), which works well with sparse gradients, and RMSProp (Tieleman & Hinton, 2012), which works well in on-line and non-stationary settings; important connections to these and other stochastic optimization methods are clarified in section 5. Some of Adam's advantages are that the magnitudes of parameter updates are invariant to rescaling of the gradient; its updates are approximately bounded by the stepsize hyperparameter; it does not require a stationary objective; it works with sparse gradients, and it naturally performs a form of step-size annealing.

¹Equal contribution. Author ordering determined by coin flip over a Google Hangout.

Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." *arXiv preprint arXiv:1412.6980* (2014).

Adam - Pytorch

pytorch.org/docs/stable/generated/torch.optim.Adam.html

Adam

```
CLASS torch.optim.Adam(params, lr=0.001, betas=(0.9, 0.999), eps=1e-08, weight_decay=0,
                        amsgrad=False, *, foreach=None, maximize=False, capturable=False, differentiable=False,
                        fused=None) [SOURCE]
```


Exercício

1. Modifique o algoritmo da aula passada, que treinava uma rede neural para classificação de vagas, para utilizar o Adam como otimizador. Compare a convergência da rede usando *learning rate* fixo (ou com adição um scheduler), versus usando Adam (com ou sem scheduler).

Anexo - Esboço da prova

$$d = \frac{1,7481 - 1,0825}{1 - 1} = 1,7481 - 1,0825 = 0,6656 \text{ g/mL}$$

d = (Pic amostra - Pic água) / (Pic amostra - Pic água)

Para cálculo da densidade relativa (d) do Alendazol a 20°C, usa-se a seguinte

Massa picnômetro com Alendazol (pic amostra): 47,7481 g
 Massa picnômetro com água (pic água): 40,8255 g
 Massa picnômetro vazio (pic vazio): 15,7261 g

Resultados:

$$\lim_{\alpha \rightarrow \infty} \frac{1}{1 - w} = 1$$

$$\lim_{\alpha \rightarrow \infty} \frac{1}{1 - w} = 1$$

Matérias: picnômetro-25 mL, água, Alendazol (suspensão).

$$p_{20} = 0,99820 \times 1,026 + 0,0017 = 1$$

p = d(água) + d + 0,0017, expressa em g/mL ou kg/L.

temperatura (t) e calculada a partir de sua densidade relativa (d) pela fórmula:
 volume a 20°C. A densidade de massa da substância (p) em uma determinada
 A densidade de massa (p) de uma substância é a razão entre sua massa por seu

$$= \frac{1}{1 - w} = \frac{1}{1 - 0,38} = 1,6129$$

PRÁTICA 2

$$\lim_{\alpha \rightarrow \infty} \frac{1}{1 - w} = 1$$

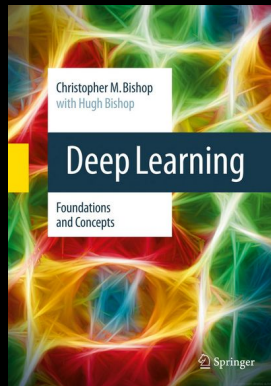
$$w^{\tau+1} = w^{\tau} - \rho \nabla E \quad \text{para } \alpha > 1$$

$$S(w, \alpha) = (1 + w + w^2 + w^3 + \dots + w^{\alpha}), \text{ para } \alpha > 1$$

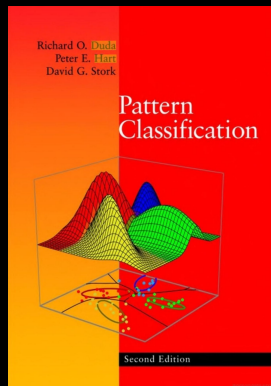


Referências

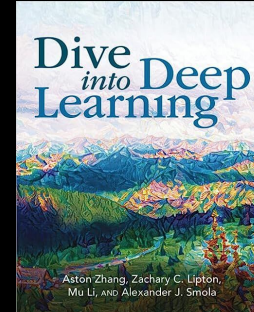
Bishop, C. M., Bishop, H. Deep Learning: Foundations and Concepts. 2024.



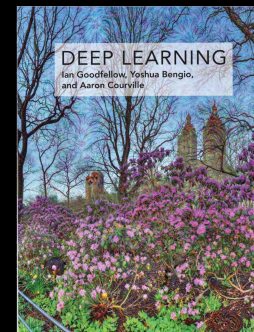
Duda, R. O., Hart, P. E., Stork, D. G. Pattern Classification. 2012.



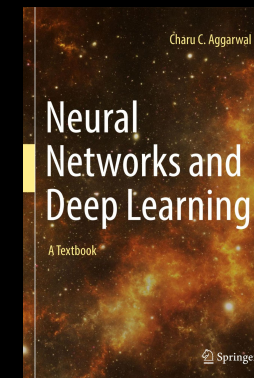
Zhang, A., Lipton, Z. C., Li, M., Smola, A. J. Dive Into Deep Learning. 2023



Goodfellow, I., Bengio, Y., Courville, A. Deep Learning. 2016.



Aggarwal. Neural Networks and Deep Learning. 2018.



Licença

Esta obra está licenciada com uma Licença [Creative Commons Atribuição 4.0 Internacional](https://creativecommons.org/licenses/by/4.0/).